

BY JULIA EVANS

10

10

10

SELECT author, COUNT(*)
FROM posts
GROUP BY author

SELECT *
FROM WHERE
ORDER BY
HAVING
LIMIT



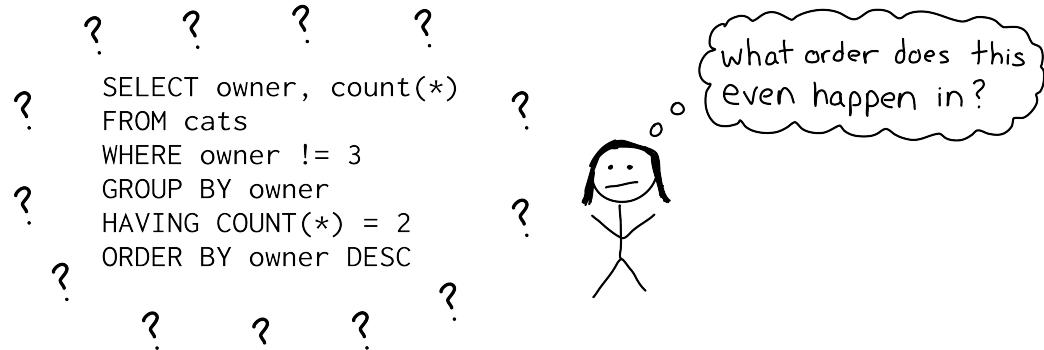
BECOME A
SELECT STAR

*wizardzines.com *

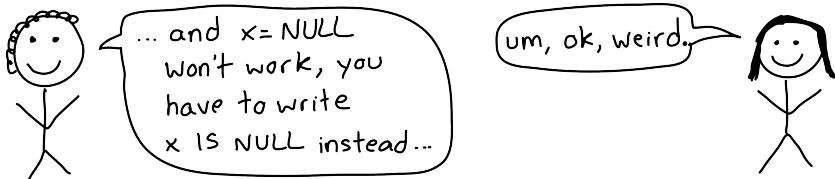
more zines at
love this?

about this zine

SQL isn't like most programming languages.



and it has quite a few weird gotchas:



but knowing it lets you use really powerful database software to answer questions FAST:



This zine will get you started with SELECT queries so you can get the answers to any question you want about your data. ❤️ You can run any query from the zine here:

<https://sql-playground.wizardzines.com>

♥ thanks for reading ♥

When you're learning it's important to experiment! So you can try out any of the queries in this zine and run your own in an SQL playground:

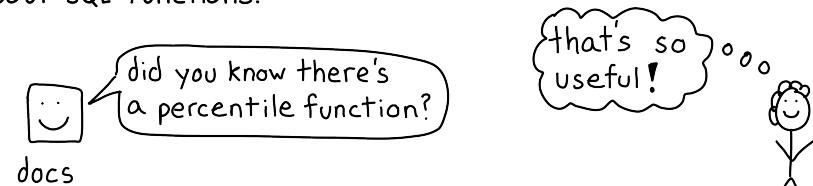
★ <https://sql-playground.wizardzines.com> ★

Here are a few more great SQL resources:

★ Use the index, Luke! (<https://use-the-index-luke.com>) is a very in-depth explanation of how to use indexes to make your queries fast.

★ There are several visualizers that will help you understand the output of an EXPLAIN. For example: <https://explain.depesz.com/> for PostgreSQL!

★ The official documentation is always GREAT for learning about SQL functions.



★ credits ★

Cover illustration: Deise Lino

Editing: Dolly Lanuza, Kamal Marhubi, Samuel Wright

Reviewers: Anton Dubrau, Arielle Evans

Table of Contents

- 4 anatomy of a SELECT query
- 5 SELECT queries start with FROM
- 6 SELECT WHERE
- 7 SELECT GROUP BY
- 8 HAVING
- 9 JOINs
- 10-11 Does this column have NULL or 0 or empty string values?
- 12 example: GROUP BY
- 13 ORDER BY + LIMIT
- 14 intro
- 15 OVER() assigns every row a window
- 16 much ado about NULL
- 17 NULL: unknown or missing
- 18 NULL surprises
- 19 handle NULLS with COALESCE
- 20 CASE
- 21 ways to count rows
- 22 subqueries
- 23 tip: single quote strings
- 24 EXPLAIN your slow queries
- 25 how indexes make your queries fast
- 26 questions to ask about your data
- 27 thanks for reading ♡

about your data

It's really easy to make incorrect assumptions about the data in a table:

• every hospital patient has a doctor right?
• why is everyone missing a doctor?
• from May 2013 3 hours later...

• doctor right?
• patient has a doctor right?
• why is everyone missing a doctor?

Some questions you might want to ask:

Does this column have NULL or 0 or empty string values?

How many different values does this column have?

Are there duplicate values in this column?

• sometimes a doctor has 2 appointments at the same time,
• that shouldn't happen ...
• why are there 213 doctor IDs with no match in the doctors table?

Does the id column in table A always have a match in table B?

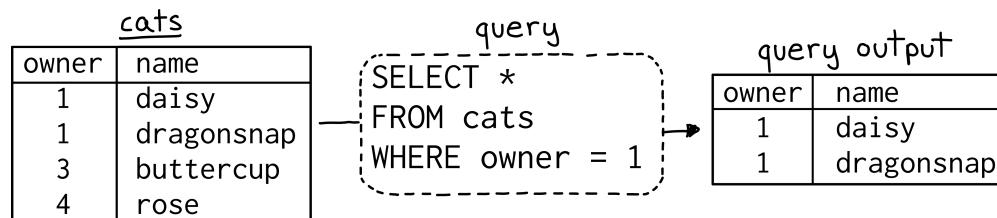
• a lot of these can also be enforced by NOT NULL or UNIQUE or FOREIGN KEY constraints on your tables.

getting started with SELECT

A SQL database contains a bunch of tables

| sales | | clients | | cats | |
|--------|------|---------|------|-------|------|
| client | item | id | name | owner | name |
| ~ | ~ | ~ | ~ | ~ | ~ |
| ~ | ~ | ~ | ~ | ~ | ~ |
| ~ | ~ | ~ | ~ | ~ | ~ |
| ~ | ~ | ~ | ~ | ~ | ~ |
| ~ | ~ | ~ | ~ | ~ | ~ |

Every SELECT query takes data from those tables and outputs a table of results.



A few basic facts to start out:

→ SELECT queries have to be written in the order:

SELECT ... FROM ... WHERE ... GROUP BY ... HAVING ...
ORDER BY ... LIMIT

→ SQL isn't case sensitive: select * from table is fine too.

This zine will use ALL CAPS for SQL Keywords like FROM.



there are other kinds of queries like INSERT/UPDATE/DELETE but this zine is just about SELECT

EXPLAIN your slow queries

Sometimes queries run slowly, and EXPLAIN can tell you why!

2 ways you can use EXPLAIN in PostgreSQL: (other databases have different syntax for this)

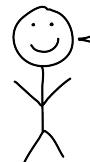
① Before running the query (EXPLAIN SELECT ... FROM ...)

This calculates a query plan but doesn't run the query.

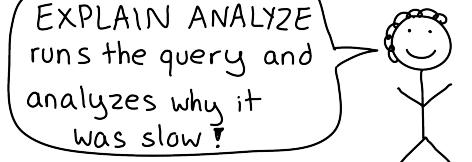


I always run EXPLAIN on a query before running it on my production database. I won't risk overloading the database with a slow query!

② After running the query (EXPLAIN ANALYZE SELECT ... FROM ...)



why is my query so slow?



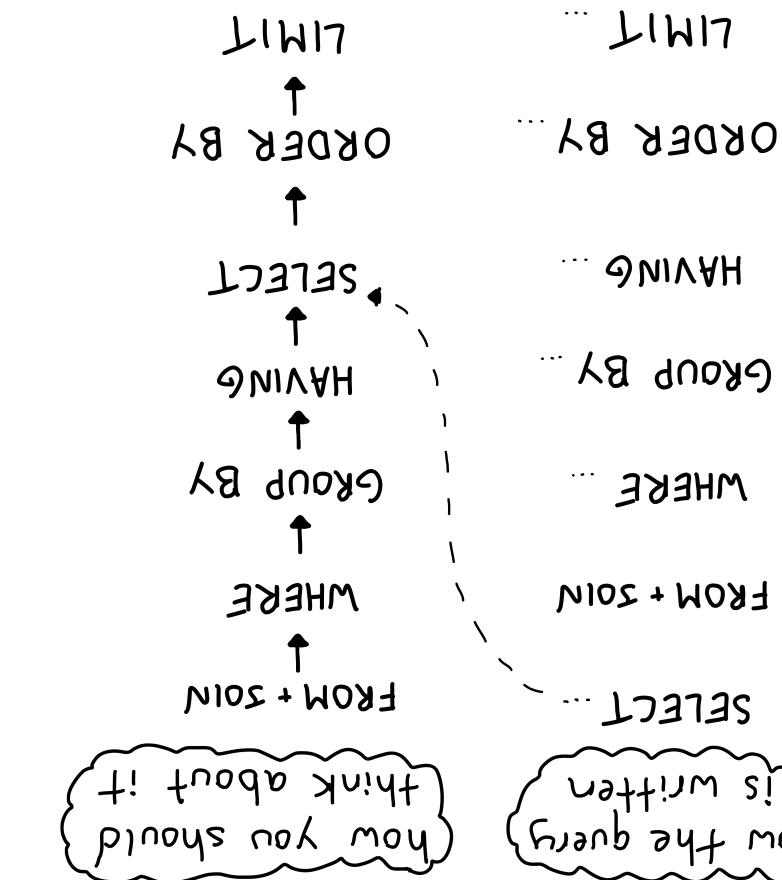
EXPLAIN ANALYZE runs the query and analyzes why it was slow!

Here are the EXPLAIN ANALYZE results from PostgreSQL for the same query run on two tables of 1,000,000 rows: one table that has an index and one that doesn't

(EXPLAIN ANALYZE SELECT * FROM users WHERE id = 1)

| unindexed table | indexed table |
|--|--|
| <p>Seq Scan on users</p> <p>Filter: (id = 1)</p> <p>Rows Removed by Filter: 999999</p> <p>Planning time: 0.185 ms</p> <p>Execution time: 179.412 ms</p> <p>"Seq Scan" means it's looking at each row (slow!)</p> | <p>Index Only Scan using users_id_idx on users</p> <p>Index Cond: (id = 1)</p> <p>Heap Fetches: 1</p> <p>Planning time: 3.411 ms</p> <p>Execution time: 0.088 ms</p> <p>the query runs 50 times faster with an index</p> |

There are a lot of optimizations.)
(In reality query execution is much more complicated than this.



The query's steps don't happen in the order they're written:

| owner | name |
|-------|------------|
| 1 | daisy |
| 1 | dragonsnap |

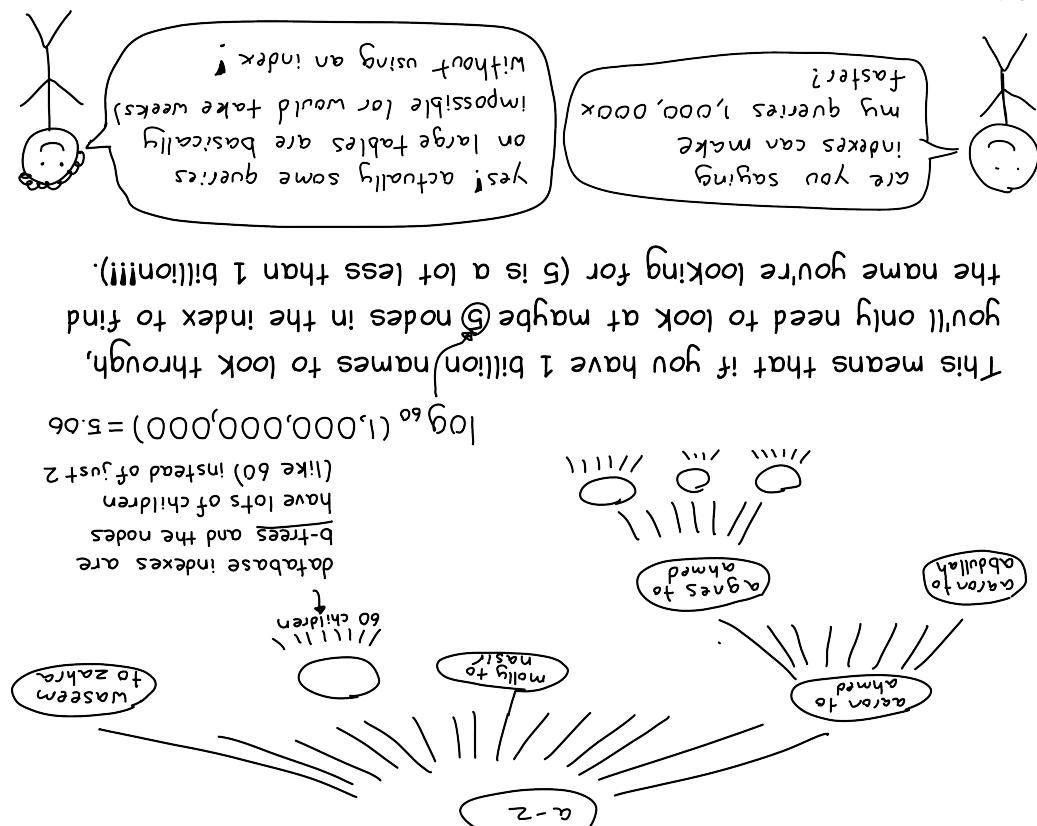
| owner | name | cat |
|-------|------------|------------|
| 1 | daisy | buttercup |
| 1 | dragonsnap | dragonSnap |
| 4 | | rosie |

WHERE owner = 1

its input, like this:

conceptually, every step (like "WHERE") of a query transforms

SELECT queries start with FROM



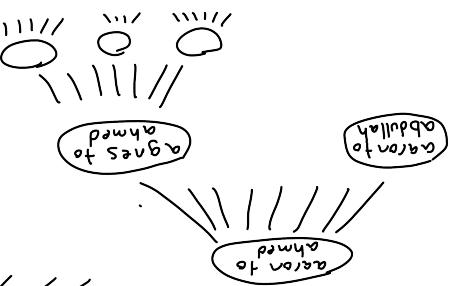
the name you're looking for (5 is a lot less than 1 billion!!!).
you'll only need to look at maybe 5 nodes in the index to find
this means that if you have 1 billion names to look through,

$$\log_{60}(1,000,000,000) = 5.66$$

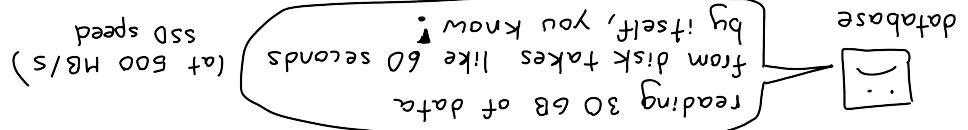
(like 60 instead of just 2
have lots of children
b-trees and the nodes
database indexes are

to zahra

60 children



Indexes are a tree structure that makes it faster to find rows.
Here's what an index on the name, column might look like.



By default, if you run SELECT * FROM cats WHERE name = 'mr darcy',
the database needs to look at every single row to find matches

your queries fast

How indexes make

SELECT

SELECT is where you pick the final columns that appear in the table the query outputs. Here's the syntax:

```
SELECT expression_1 [AS alias],  
      expression_2 [AS alias2],  
  FROM ...
```

Some useful things you can do in a SELECT:

★ Combine many columns with SQL expressions

A few examples:

```
CONCAT(first_name, ' ', last_name)  
DATE_TRUNC('month', created) ← this is PostgreSQL syntax for  
                           rounding a date, other SQL  
                           dialects have different syntax
```

★ Alias an expression with AS

first_name || ' ' || last_name is a mouthful! If you alias an expression with AS, you can use the alias elsewhere in the query to refer to that expression.

```
SELECT first_name || ' ' || last_name AS full_name  
FROM people  
ORDER BY full_name DESC  
        ↑ refers to first_name || ' ' || last_name
```

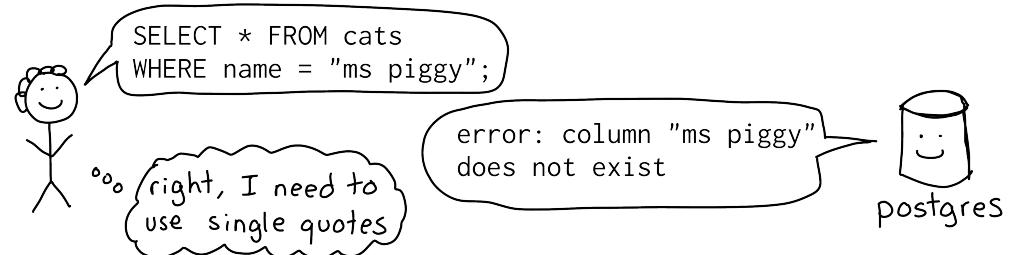
★ Select all columns with SELECT *

When I'm starting to figure out a query, I'll often write something like

```
SELECT * FROM some_table LIMIT 10  
just to quickly see what the columns in the table look like.
```

tip: single quote strings

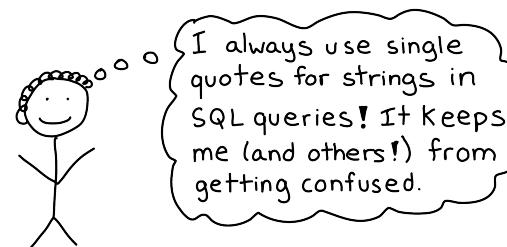
In some SQL implementations (like PostgreSQL), if you double quote a string it'll interpret it as a column name:



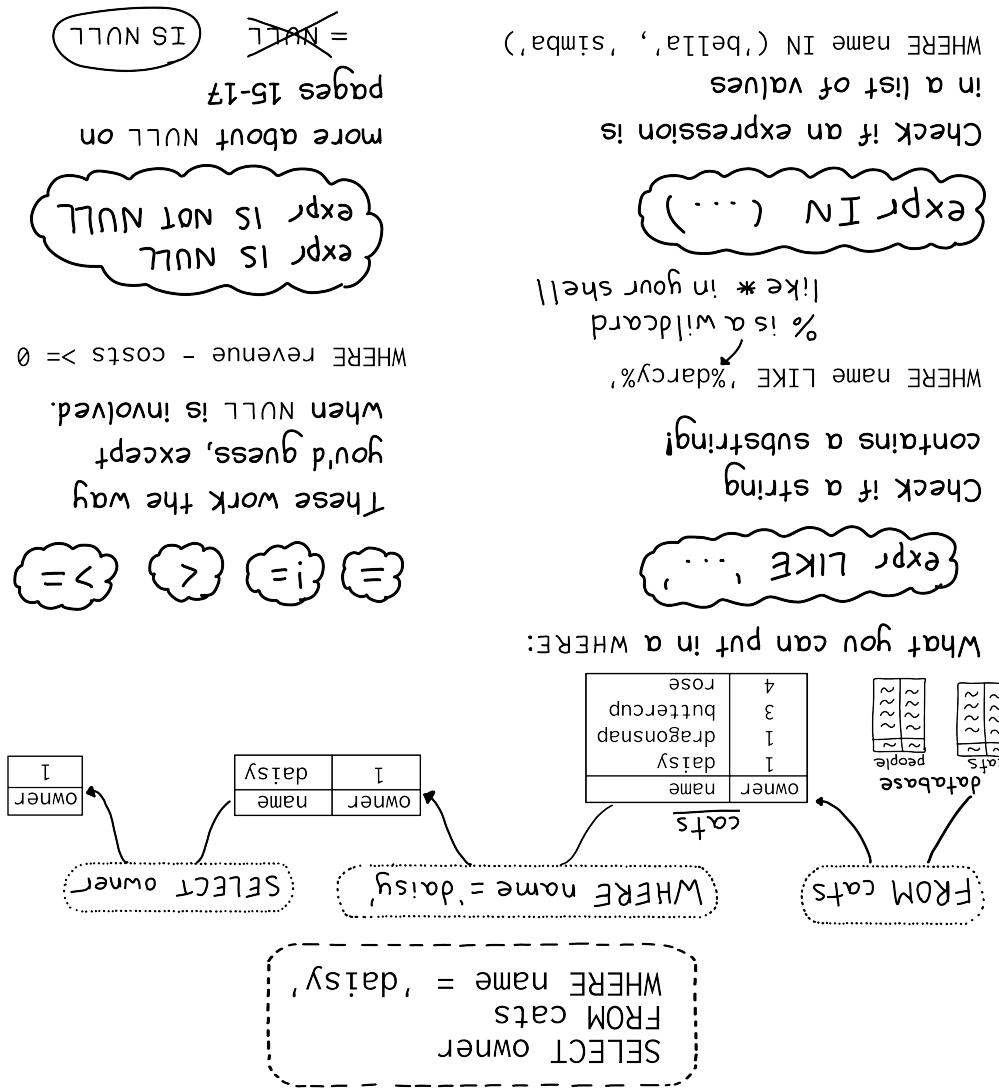
Here's a table explaining what different quotes mean in different SQL databases. "Identifier" means a column name or table name.

Sometimes table names have special characters like spaces in them so it's useful to be able to quote them.

| | single quotes 'ms piggy' | double quotes "ms piggy" | backticks '`ms piggy`' |
|------------|-----------------------------|-----------------------------|---------------------------|
| MySQL | string | string or identifier | identifier |
| PostgreSQL | string | identifier | invalid |
| SQLite | string | string or identifier | identifier |
| SQL server | string | string or identifier | invalid |



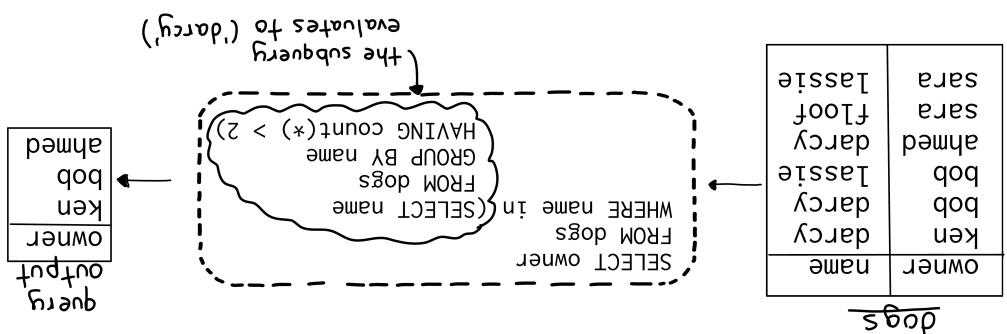
7
WHERE filters the table you start with. For example, let's break down this query that finds all owners with cats named "daisy".
SELECT owner FROM cats WHERE name = 'daisy'
popular names: ("boring" owners):
For example, this query finds owners who have named their dogs
the subquery evaluates to (darcy)
SELECT owner WHERE name IN (SELECT name FROM dogs WHERE count(*) > 2)
SELECT owner FROM dogs GROUP BY name HAVING count(*) > 2
SELECT owner FROM popular_dog_names AS subquery WHERE name IN (SELECT name FROM dogs WHERE count(*) > 2)
SELECT owner FROM popular_dog_names JOIN popular_dog_names ON dogs.name = popular_dog_names.name
SELECT owner FROM dogs GROUP BY name HAVING count(*) > 2
WITH popular_dog_names AS (SELECT name FROM dogs GROUP BY name HAVING count(*) > 2)
SELECT owner FROM owner WHERE expr IS NOT NULL
expr IS NOT NULL
expr IS NULL
more about NULL on pages 15-17
WHERE name IN ('belila', 'simba')
in a list of values
Check if an expression is
expr IN (...)
like * in your shell
% is a wildcard
WHERE name LIKE '%darcy%'
contains a substring!
Check if a string
expr LIKE '...'
These work the way
you'd guess, except
when NULL is involved.
WHERE revenue - costs >= 0
WHERE name LIKE '...'
expr IS NULL
expr IS NOT NULL
IS NULL
= NOT
AND OR NOT
You can AND together as many conditions as you want
If I'm using lots of ANDS I like to write them like this
in the parentheses put all the ORs
AND (....)
AND (....)
AND (....)
AND (....)



Some questions can't be answered with one simple SQL query.

Subqueries

For example, this query finds owners who have named their dogs
the subquery evaluates to (darcy)
SELECT owner FROM owner WHERE name IN (SELECT name FROM dogs WHERE count(*) > 2)
SELECT owner FROM dogs GROUP BY name HAVING count(*) > 2
WITH popular_dog_names AS (SELECT name FROM dogs GROUP BY name HAVING count(*) > 2)
SELECT owner FROM owner WHERE expr IS NOT NULL
expr IS NOT NULL
expr IS NULL
more about NULL on pages 15-17
WHERE name IN ('belila', 'simba')
in a list of values
Check if an expression is
expr IN (...)
like * in your shell
% is a wildcard
WHERE name LIKE '%darcy%'
contains a substring!
Check if a string
expr LIKE '...'
These work the way
you'd guess, except
when NULL is involved.
WHERE revenue - costs >= 0
WHERE name LIKE '...'
expr IS NULL
expr IS NOT NULL
IS NULL
= NOT
AND OR NOT
You can AND together as many conditions as you want
If I'm using lots of ANDS I like to write them like this
in the parentheses put all the ORs
AND (....)
AND (....)
AND (....)
AND (....)

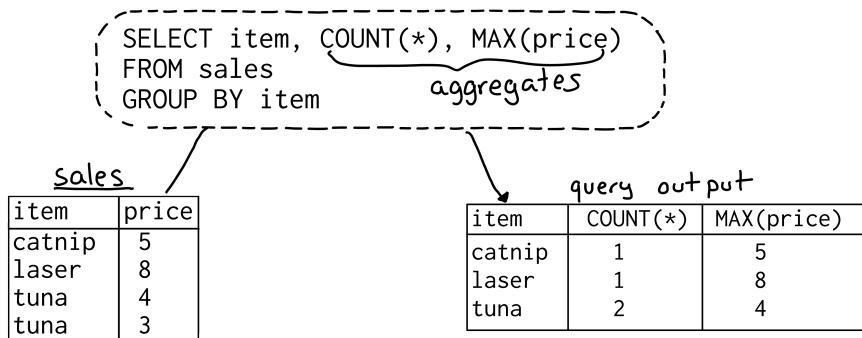


For example, this query finds owners who have named their dogs
the subquery evaluates to (darcy)
SELECT owner FROM owner WHERE name IN (SELECT name FROM dogs WHERE count(*) > 2)
SELECT owner FROM dogs GROUP BY name HAVING count(*) > 2
WITH popular_dog_names AS (SELECT name FROM dogs GROUP BY name HAVING count(*) > 2)
SELECT owner FROM owner WHERE expr IS NOT NULL
expr IS NOT NULL
expr IS NULL
more about NULL on pages 15-17
WHERE name IN ('belila', 'simba')
in a list of values
Check if an expression is
expr IN (...)
like * in your shell
% is a wildcard
WHERE name LIKE '%darcy%'
contains a substring!
Check if a string
expr LIKE '...'
These work the way
you'd guess, except
when NULL is involved.
WHERE revenue - costs >= 0
WHERE name LIKE '...'
expr IS NULL
expr IS NOT NULL
IS NULL
= NOT
AND OR NOT
You can AND together as many conditions as you want
If I'm using lots of ANDS I like to write them like this
in the parentheses put all the ORs
AND (....)
AND (....)
AND (....)
AND (....)

GROUP BY
FROM (<subquery or CTE>)
SELECT
WHERE name IN (<subquery>)

GROUP BY

GROUP BY combines multiple rows into one row. Here's how it works for this table & query:



① Split the table into groups for each value that you grouped by:

| item='catnip' | item='laser' | item='tuna' | | | | | | | | | | | | | | |
|---|--------------|-------------|--------|---|--|------|-------|-------|---|---|------|-------|------|---|------|---|
| <table border="1"> <thead> <tr> <th>item</th> <th>price</th> </tr> </thead> <tbody> <tr> <td>catnip</td> <td>5</td> </tr> </tbody> </table> | item | price | catnip | 5 | <table border="1"> <thead> <tr> <th>item</th> <th>price</th> </tr> </thead> <tbody> <tr> <td>laser</td> <td>8</td> </tr> </tbody> </table> | item | price | laser | 8 | <table border="1"> <thead> <tr> <th>item</th> <th>price</th> </tr> </thead> <tbody> <tr> <td>tuna</td> <td>4</td> </tr> <tr> <td>tuna</td> <td>3</td> </tr> </tbody> </table> | item | price | tuna | 4 | tuna | 3 |
| item | price | | | | | | | | | | | | | | | |
| catnip | 5 | | | | | | | | | | | | | | | |
| item | price | | | | | | | | | | | | | | | |
| laser | 8 | | | | | | | | | | | | | | | |
| item | price | | | | | | | | | | | | | | | |
| tuna | 4 | | | | | | | | | | | | | | | |
| tuna | 3 | | | | | | | | | | | | | | | |

② Calculate the aggregates from the query for each group:

| <table border="1"> <thead> <tr> <th>item</th> <th>price</th> </tr> </thead> <tbody> <tr> <td>catnip</td> <td>5</td> </tr> </tbody> </table> | item | price | catnip | 5 | <table border="1"> <thead> <tr> <th>item</th> <th>price</th> </tr> </thead> <tbody> <tr> <td>laser</td> <td>8</td> </tr> </tbody> </table> | item | price | laser | 8 | <table border="1"> <thead> <tr> <th>item</th> <th>price</th> </tr> </thead> <tbody> <tr> <td>tuna</td> <td>4</td> </tr> <tr> <td>tuna</td> <td>3</td> </tr> </tbody> </table> | item | price | tuna | 4 | tuna | 3 |
|---|----------------|----------------|--------|---|--|------|-------|-------|---|---|------|-------|------|---|------|---|
| item | price | | | | | | | | | | | | | | | |
| catnip | 5 | | | | | | | | | | | | | | | |
| item | price | | | | | | | | | | | | | | | |
| laser | 8 | | | | | | | | | | | | | | | |
| item | price | | | | | | | | | | | | | | | |
| tuna | 4 | | | | | | | | | | | | | | | |
| tuna | 3 | | | | | | | | | | | | | | | |
| COUNT(*) = 1 | COUNT(*) = 1 | COUNT(*) = 2 | | | | | | | | | | | | | | |
| MAX(price) = 5 | MAX(price) = 8 | MAX(price) = 4 | | | | | | | | | | | | | | |

③ Create a result set with 1 row for each group

| item | COUNT(*) | MAX(price) |
|--------|----------|------------|
| catnip | 1 | 5 |
| laser | 1 | 8 |
| tuna | 2 | 4 |

ways to count rows

Here are three ways to count rows:

① COUNT(*): count all rows

This counts every row, regardless of the values in the row.
Often used with a GROUP BY to get common values, like in this "most popular names" query:

```
SELECT first_name, COUNT(*)
FROM people
GROUP BY first_name
ORDER BY COUNT(*) DESC
LIMIT 50
```

② COUNT(DISTINCT column): get the number of distinct values

Really useful when a column has duplicate values.
For example, this query finds out how many species every plant genus has:

"GROUP BY 1"
means group by
the first expression
in the SELECT

```
SELECT genus, COUNT(DISTINCT species)
FROM plants
GROUP BY 1
ORDER BY 2 DESC
```

③ SUM(CASE WHEN expression THEN 1 ELSE 0 END)

This trick using SUM and CASE lets you count how many cats vs dogs vs other animals each owner has:

I like to put commas at the start for big queries

```
SELECT owner
, SUM(CASE WHEN type = 'dog' then 1 else 0 end) AS num_dogs
, SUM(CASE WHEN type = 'cat' then 1 else 0 end) AS num_cats
, SUM(CASE WHEN type NOT IN ('dog', 'cat') then 1 else 0
end) AS num_other
FROM pets GROUP BY owner
```

| owner | type |
|-------|----------|
| 1 | dog |
| 1 | cat |
| 2 | dog |
| 2 | parakeet |

| owner | num_dogs | num_cats | num_other |
|-------|----------|----------|-----------|
| 1 | 1 | 1 | 0 |
| 2 | 1 | 0 | 1 |

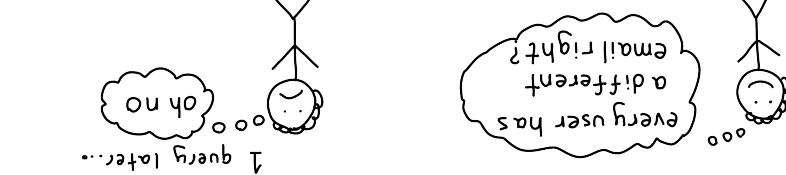
HAVING uses HAVING to find all emails that are shared by more than one user:

This query uses HAVING to find all emails that are shared by more than one user:

SELECT email, COUNT(*)
FROM users
GROUP BY email
HAVING COUNT(*) > 1

| id | email | COUNT(*) |
|----|-----------------|----------|
| 1 | asdfefake.com | 1 |
| 2 | asdfefake.com | 1 |
| 3 | boblebulder.com | 1 |

query output



HAVING is like WHERE, but with 1 difference: HAVING filters rows AFTER grouping and WHERE filters rows BEFORE grouping.

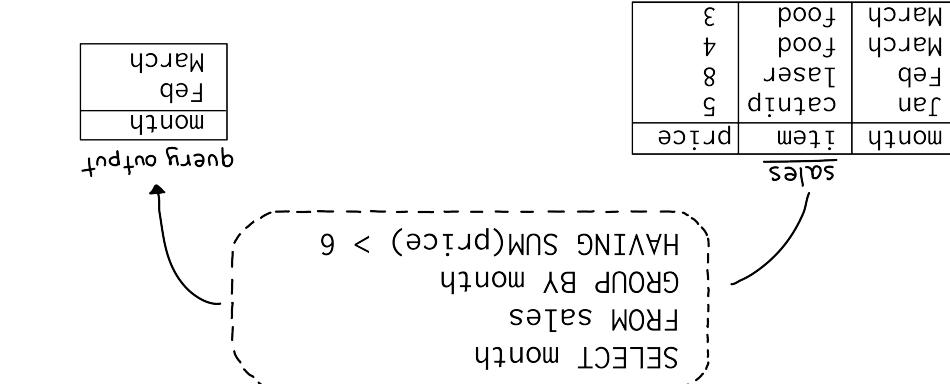
Because of this, you can use aggregates (like COUNT(*)) in a HAVING clause but not with WHERE.

Here's another HAVING example that finds months with more than \$6.00 in income:

SELECT month FROM sales GROUP BY month HAVING SUM(price) > 6

| month | item | price | sales | category | food | March |
|-------|-------|-------|-------|----------|------|-------|
| Jan | laser | 8 | | | 4 | |
| Feb | laser | 5 | | | 3 | |
| Mar | laser | 8 | | | 4 | |

query output



HAVING

Here's how to categorize people into age ranges!

*** Example ***

| first_name | age | age_range |
|------------|-----|-----------|
| ahmed | 5 | child |
| marie | 17 | adult |
| akira | 60 | teenager |
| pablo | 15 | adult |

people

SELECT first_name, age, age_range
FROM people
WHEN age < 13 THEN 'child'
WHEN age < 20 THEN 'adult'
ELSE 'adult', END AS age_range
RETURNS first_name
condition matches
SELECT first_name, age, age_range
FROM people
WHEN age < 13 THEN 'child',
WHEN age < 20 THEN 'teenager',
ELSE 'adult', END AS age_range
RETURNS first_name
condition matches

CASE is how to write an if statement in SQL. Here's the syntax:

```

CASE
  WHEN <Condition> THEN <result>
  WHEN <Condition> THEN <result>
  ...
  ELSE <result>
END

```

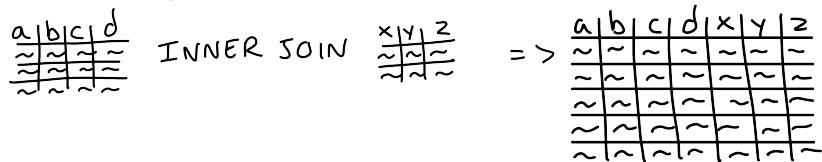


Often I want to categorize by something that isn't a column:

CASE

* my rules for simple JOINS *

Joins in SQL let you take 2 tables and combine them into one.



Joins can get really complicated, so we'll start with the simplest way to join. Here are the rules I use for 90% of my joins:

Rule 1: only use LEFT JOIN and INNER JOIN

There are other kinds of joins (RIGHT JOIN, CROSS JOIN, FULL OUTER JOIN), but 95% of the time I only use LEFT JOIN and INNER JOIN.

Rule 2: refer to columns as table_name.column_name

You can leave out the table name if there's just one column with that name, but it can get confusing

Rule 3: Only include 1 condition in your join

Here's the syntax for a LEFT JOIN:

table1 LEFT JOIN table2 ON <any boolean condition>

I usually stick to a very simple condition, like this:

```
table1 LEFT JOIN table2
  ON table1.some_column = table2.other_column
```

Rule 4: One of the joined columns should have unique values

If neither of the columns is unique, you'll get strange results like this:

| owners_bad | | cats_bad | | owners_bad INNER JOIN cats_bad ON owners_bad.name = cats_bad.owner | | |
|------------|-----|----------|------------|---|------------|-----|
| name | age | owner | name | owner | name | age |
| maher | 16 | maher | daisy | maher | daisy | 16 |
| maher | 32 | maher | dragonsnap | maher | dragonsnap | 16 |
| rishi | 21 | rishi | buttercup | rishi | daisy | 32 |
| | | | | | maher | 32 |
| | | | | | dragonsnap | 21 |

INNER JOIN

(these are "bad" versions of the `owners` and `cats` tables that don't JOIN well)

handle NULLs with COALESCE

COALESCE is a function that returns the first argument you give it that isn't NULL

```
COALESCE(NULL, 1, 2) => 1
COALESCE(NULL, NULL, NULL) => NULL
COALESCE(4, NULL, 2) => 4
```

2 ways you might want to use COALESCE in practice:

① Set a default value

In this table, a NULL discount means there's no discount, so we use COALESCE to set the default to 0:

products

| name | price | discount |
|--------|-------|----------|
| orange | 200 | NULL |
| apple | 100 | 23 |
| lemon | 150 | NULL |

query output

| name | net_price |
|--------|-----------|
| orange | 200 |
| apple | 77 |
| lemon | 150 |

② Use data from 2 (or more!) different columns

This query gets the best guess at a customer's state:

mailing address, most accurate

if not, try billing address

as a last resort, use their IP address

addresses

| customer | mailing_state | billing_state | ip_address_state |
|----------|---------------|---------------|------------------|
| 1 | Bihar | Bihar | Bihar |
| 2 | NULL | Kerala | Kerala |
| 3 | NULL | NULL | Punjab |
| 4 | Gujarat | Punjab | Gujarat |

state

| state |
|---------|
| Bihar |
| Kerala |
| Punjab |
| Gujarat |

- Here are examples of how INNER JOIN and LEFT JOIN work:
- INNER JOIN only includes rows that match the ON condition.
 - This query combines the cats and owners tables:
 - SELECT * FROM owners INNER JOIN cats ON cats.owner = owners.id
 - These 2 rows don't have matches so they're not in the output.
 - LEFT JOIN includes every row from the left table (owners in this example), even if it's not in the right table. Rows not in the right table will be set to NULL.
 - SELECT * FROM owners LEFT JOIN cats ON cats.owner = owners.id
 - These 2 rows don't have matches so they're not in the output.

This is a classic example of a join that follows my 4 guidelines from the previous page:

Rishi has no cats

| cats.name | owners.name |
|------------|-------------|
| rose | buttercup |
| chandara | daisy |
| maher | rishi |
| dragonsnap | dragnsnap |
| | NULL |

| owners.id | owners.name |
|-----------|-------------|
| 1 | daisy |
| 2 | maher |
| 3 | rishi |
| 4 | chandara |

1) it's an INNER JOIN / LEFT JOIN

2) it includes the table name in cats.owner and owners.id

3) the condition ON cats.owner = owners.id is simple

4) it joins on a unique column (the id column in the owners table)

More operations with NULL which might be surprising: NULL isn't even equal to itself!

| | |
|-------------|---------|
| 2 + NULL | => NULL |
| NULL * 10 | => NULL |
| NULL = NULL | => NULL |
| NULL = NULL | => NULL |

CONCAT('hi', NULL) => NULL

2 = NULL => NULL

2 != NULL => NULL

more surprising truths

To match NULLs as well, I'll often write something like WHERE name != 'bette', OR name IS NULL instead.

| | | |
|------|-------|-------|
| name | owner | ahmed |
| name | owner | nemo |
| name | owner | bob |

SURPRISE! name != 'bette', doesn't match NULLs

| | |
|------|-------------|
| name | owner |
| name | NULL |
| name | IS NOT NULL |

works

| | |
|------|---------|
| name | owner |
| name | NULL |
| name | IS NULL |

You need to use != IS NULL instead.

| | |
|------|-------------|
| name | owner |
| name | NULL |
| name | IS NOT NULL |

SURPRISE! x=NULL doesn't work

| | |
|------|---------|
| name | owner |
| name | NULL |
| name | IS NULL |

that are surprising at first:

x != NULL are never true for any x). This results in 2 behaviours

NULL isn't equal (or not equal) to anything in SQL (x = NULL and

NULL surprises

INNER JOIN and LEFT JOIN

example : LEFT JOIN + GROUP BY

This query counts how many items every client bought
(including clients who didn't buy anything):

```
SELECT name, COUNT(item) AS items_bought
FROM owners LEFT JOIN sales
    ON owners.id = sales.client
GROUP BY name
ORDER BY items_bought DESC
```

FROM owners LEFT JOIN sales...

| owners | |
|--------|---------|
| id | name |
| 1 | maher |
| 2 | rishi |
| 3 | chandra |

| sales | |
|--------|--------|
| item | client |
| catnip | 1 |
| laser | 1 |
| tuna | 1 |
| tuna | 2 |

ON owners.id = sales.client

GROUP BY name

| id | name | item |
|----|---------|--------|
| 1 | maher | catnip |
| 1 | maher | laser |
| 1 | maher | tuna |
| 2 | rishi | tuna |
| 3 | chandra | NULL |

| id | name | item |
|----|---------|--------|
| 1 | maher | catnip |
| 1 | maher | laser |
| 1 | maher | tuna |
| 2 | rishi | tuna |
| 3 | chandra | NULL |

SELECT name, COUNT(item)
AS items_bought

ORDER BY
items_bought DESC

| name | items_bought |
|---------|--------------|
| rishi | 1 |
| chandra | 0 |
| maher | 3 |

COUNT(item)
doesn't count
NULLs

| name | items_bought |
|---------|--------------|
| maher | 3 |
| rishi | 1 |
| chandra | 0 |

NULL: unknown or missing

NULL is a special state in SQL. It's very commonly used as a placeholder for missing data ("we don't know her address!")

What NULL means exactly depends on your data. For example, it's really important to know if allergies IS NULL means

→ "no allergies" or

→ "we don't know if she has allergies or not"

NULL "should" mean
"unknown" but it doesn't
always



it would be way easier
if NULL always meant
the same thing but it
really depends on
your data!

★ where NULLs come from ★

→ There were already NULL values in the table

→ The window function LAG() can return NULL

→ You did a LEFT JOIN and some of the rows on the left didn't have a match for the ON condition

ooh, not every cat has
an owner so sometimes
the owner name is NULL

★ ways to handle NULLs ★

→ Leave them in!

I'd rather see a NULL
and know there's missing data
than get misleading results

→ Filter them out!

... WHERE first_name IS NOT NULL ...

→ Use COALESCE or CASE to add a default value

| | |
|-------|-----------|
| owner | name |
| 4 | rose |
| 3 | daisy |
| 1 | buttercup |

| | |
|-------|------------|
| owner | name |
| 1 | daisy |
| 3 | dragonsnap |
| 4 | buttercup |

SELECT * FROM cats
ORDER BY LENGTH(name) ASC
LIMIT 2

For example, this is the same as the previous query, but it limits to only the 2 cats with the shortest names:

LIMIT [integer]

LIMIT lets you limit the number of rows output.

The syntax is:

| | |
|-------|-----------|
| owner | name |
| 4 | rose |
| 3 | daisy |
| 1 | buttercup |

| | |
|-------|------------|
| owner | name |
| 1 | daisy |
| 3 | dragonsnap |
| 4 | buttercup |

SELECT * FROM cats
ORDER BY LENGTH(name) ASC
(shortest first):

For example, this query sorts cats by the length of their name

ORDER BY [expression] DESC

ORDER BY lets you sort by anything you want!
ASC → ascending
DESC → descending
stands for

The syntax is:

ORDER BY lets you sort by anything you want!

ORDER BY and LIMIT happen at the end and affect the final output of the query.

ORDER BY AND LIMIT

| | | |
|-------|---------------------------|-----------------|
| event | hour | time_since_Last |
| 1 | NULL | LAG() |
| 3 | is NULL for the first row | in the window |
| 4 | diaper | feeding |
| 5 | diaper | diaper |

④ SELECT + type, hour,
hour - LAG(hour)

| | | |
|-------|---------------------------|-----------------|
| event | hour | time_since_Last |
| 1 | NULL | LAG() |
| 3 | is NULL for the first row | in the window |
| 4 | diaper | feeding |
| 5 | diaper | diaper |

⑤ ORDER BY hour ASC

| | | |
|-------|------|-----------------|
| event | hour | time_since_Last |
| 1 | 1 | LAG() |
| 3 | 4 | feeding |
| 4 | 7 | diaper |
| 5 | 5 | diaper |

⑥ ORDER BY event ORDER BY hour ASC

| | | |
|-------|------|-----------------|
| event | hour | time_since_Last |
| 1 | 1 | coug |
| 3 | 3 | feeding |
| 4 | 4 | diaper |
| 5 | 5 | diaper |

⑦ FROM baby-log

| | | |
|-------|------|-----------------|
| event | hour | time_since_Last |
| 1 | 1 | coug |
| 3 | 3 | feeding |
| 4 | 4 | diaper |
| 5 | 5 | diaper |

⑧ SELECT event, hour,
hour - LAG(hour) OVER(PARTITION BY event ORDER BY hour ASC)

This query finds the time since a baby's last feeding/diaper change.

example: get the time
between baby feedings

refer to other rows with ★ window functions ★

Let's talk about an ~~advanced~~ SQL feature: window functions!

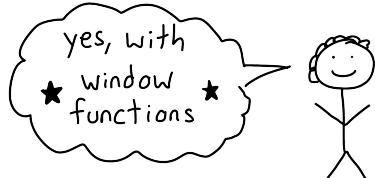
Normally SQL expressions only let you refer to information in a single row.

2 columns from the same row

```
SELECT CONCAT(firstname, ' ', lastname) as full_name
```



can I refer to other rows though? Like subtract the value in the previous row?



Window functions are SQL expressions that let you reference values in other rows. The syntax (explained on the next page!) is:

[expression] OVER ([window definition])

Example: use LAG() to find how long since the last sale

| sales | | item | day - LAG(day) OVER (ORDER BY day) |
|--------|-----|--------|------------------------------------|
| item | day | | |
| catnip | 2 | catnip | NULL |
| laser | 40 | laser | 38 |
| tuna | 70 | tuna | 30 |
| tuna | 72 | tuna | 2 |

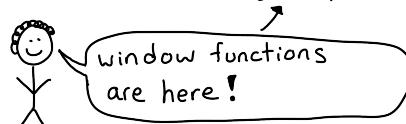
query output

| item | day - LAG(day) OVER (ORDER BY day) |
|--------|------------------------------------|
| catnip | NULL |
| laser | 38 |
| tuna | 30 |
| tuna | 2 |

2-NULL
40-2
70-40
72-70

They're part of SELECT, so they happen after HAVING:

FROM + JOIN → WHERE → GROUP BY → HAVING → SELECT → ORDER BY → LIMIT



OVER() assigns every row a window

A "window" is a set of rows.

| name | class | grade |
|-------|-------|-------|
| juan | 1 | 93 |
| lucia | 1 | 98 |

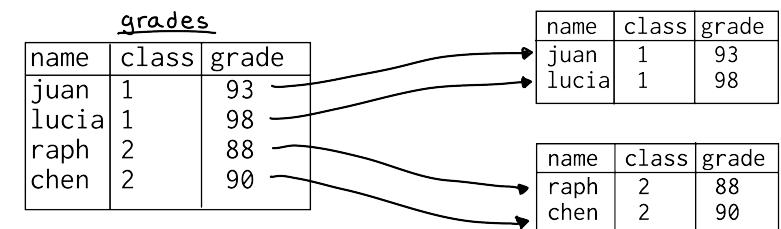
a window !

A window can be as big as the whole table (an empty OVER() is the whole table!) or as small as just one row.

OVER() is confusing at first, so here's an example! Let's run this query that ranks students in each class by grade:

```
SELECT name, class, grade,
       ROW_NUMBER() OVER (PARTITION BY class
                           ORDER BY grade DESC)
              AS rank_in_class
  FROM grades
```

Step 1: Assign every row a window. OVER(PARTITION BY class) means that there are 2 windows: one each for class 1 and 2



Step 2: Run the function. We need to run ROW_NUMBER() to find each row's rank in its window:

query output

| name | class | grade | rank_in_class |
|-------|-------|-------|---------------|
| juan | 1 | 93 | 2 |
| lucia | 1 | 98 | 1 |
| raph | 2 | 88 | 2 |
| chen | 2 | 90 | 1 |